# Relative Lempel-Ziv Compression of Suffix Arrays

Simon J. Puglisi and Bella Zhukova

# Problem (Pattern Matching)

Find all occurrences of a pattern P in a text T

Popular solution: find the interval of the suffix array (SA) that contains them

- Binary search using SA and text, or

- Backward search on the Burrows-Wheeler Transform of T (FM-index)

- Lots of compressed versions of the SA

  – Problem then becomes: how do we decompress the interval's contents?

# Previous work decoding intervals

- r-index (Gagie et al., SODA 2018)

  - recent, very fast, very small — a huge leap forward in compressed indexing

- CDAWG: succinct acyclic word graph (Belazzougui et al., CPM 2015)

  - faster than r-index on some data

  - current implementations only work for DNA

# Our contribution

A compressed SA that is bigger than the r-index, smaller than the CDAWG, and much much faster in practice than both

# Our interest in the problem

Our recent algorithms for the variable-length gap pattern matching problem (SOFSEM 2020) make scans of intervals of uncompressed SA

How to compress a SA so that decompression of random intervals would be fast?

# Core idea

| $SA$ | | | | | 30 | 25 | 20 | 15 | 10 | 5 | | | | | | | | | | | 29 | 24 | 19 | 14 | 9 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

repetitions that are off by 1 (Mäkinen, CPM 2000)

differences will turn into actual repetitions (González, Navarro, CPM 2007)

# Overview

Compression:

1. form differentially encoded $SA^{diff}$ from $SA$

2. form reference $R$ by selecting substrings from $SA^{diff}$

3. use Relative Lempel-Ziv (RLZ) to parse $SA^{diff}$ relative to $R$

4. output reference R plus set of phrases (pointers into R)

Decompression requires:

1. predecessor data structure containing phrase starting positions (in order to find the phrase covering the start of an interval)

2. absolute $SA$ value for a starting position of the phrase

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |

$x$          $x'$

$SA[x, y]$ preceded by symbol $c \Rightarrow$
$$\exists SA[x', x' + (y - x)] :$$
$$\forall i \in [0, y - x] \quad SA[x + i] = SA[x' + i] + 1$$

(González and Navarro, CPM 2007)

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$x'$          $x$

$$SA[x, y] \text{ preceded by symbol } c \Rightarrow$$
$$\exists SA[x', x' + (y - x)] :$$
$$\forall i \in [0, y - x] \quad SA[x + i] = SA[x' + i] + 1$$

(González and Navarro, CPM 2007)

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

$$i \in [1, n-1] \qquad SA^{diff}[i] = SA[i] - SA[i-1] + n$$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 | | | | | | | | | | | | | | | | | | | | |

phrases: $P$ $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

phrases

$P$      $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | \$ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** |  |  |  |  |  |  |  | 0 |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$      $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | \$ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference: | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |

phrases: **31** | | | | | | | | 0 | | | | | | | |

$P$ ⎵  $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | | | | | | | 0 | 1 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\qquad\qquad\qquad}_{P}$ $\underbrace{\qquad\qquad\qquad\qquad\qquad}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference: 27 27 27 27 27 55 27 27 27 27 27 27

phrases: **31** **30** | | | | | | | 0 | 1 | | | | | | 

$P$  $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | | | | | | | 0 | 1 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\hspace{3cm}}_{P} \qquad \underbrace{\hspace{4cm}}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | | | | | | | 0 | 1 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$ $\qquad$ $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference: 27 27 27 27 27 55 27 27 27 27 27 27

phrases: **31** **30** ... 0 1 ...

$P$ $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference: 27 27 27 27 27 55 27 27 27 27 27 27

phrases: **31** **30** _ _ _ _ _ _ | 0 1 _ _ _ _ _ _

$P$   $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 | | | | | | | | | | | | | | | | | | | | |
| phrases | **31** | **30** | 6 | **28** | | | | | 0 | 1 | 2 | 8 | | | | | | | | | | | | | | | | | | | | |

$P$       $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 | | | | | | | | | | | | | | | | | | | | |
| phrases | **31** | **30** | 6 | **28** | | | | | 0 | 1 | 2 | 8 | | | | | | | | | | | | | | | | | | | | |

$P$    $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | 6 | **28** | | | | | 0 | 1 | 2 | 8 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\hphantom{}}_{P}$ $\underbrace{\hphantom{}}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference

| 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|

phrases

| **31** | **30** | 6 | **28** | | | | | 0 | 1 | 2 | 8 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\qquad\qquad}_{P}$  $\underbrace{\qquad\qquad}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | 6 | **28** | 0 | **29** | | | 0 | 1 | 2 | 8 | 9 | 20 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$ spans the first group; $S$ spans the second group.

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference: 27 27 27 27 27 55 27 27 27 27 27 27

phrases: **31** **30** 6 **28** 0 **29** | 0 1 2 8 9 20

$\underbrace{\qquad\qquad}_{P}$ $\underbrace{\qquad\qquad}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference

| 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|

phrases

| **31** | **30** | 6 | **28** | 0 | **29** | | | 0 | 1 | 2 | 8 | 9 | 20 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\qquad\qquad}_{P} \qquad \underbrace{\qquad\qquad}_{S}$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | 6 | **28** | 0 | **29** | | | 0 | 1 | 2 | 8 | 9 | 20 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$ spans the first group; $S$ spans the second group.

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | 6 | **28** | 0 | **29** | | | 0 | 1 | 2 | 8 | 9 | 20 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$     $S$

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |
| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |

reference:

| 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|

phrases:

| **31** | **30** | 6 | **28** | 0 | **29** | 0 | - | 0 | 1 | 2 | 8 | 9 | 20 | 21 | 32 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

(The first eight entries form group $P$; the last eight entries form group $S$.)

# Example

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | c | t | a | g | a | $ |

| $SA$ | 31 | 30 | 25 | 20 | 15 | 10 | 5 | 0 | 28 | 23 | 18 | 13 | 8 | 3 | 26 | 21 | 16 | 11 | 6 | 1 | 29 | 24 | 19 | 14 | 9 | 4 | 27 | 22 | 17 | 12 | 7 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| $SA^{diff}$ | 31 | 31 | 27 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 60 | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| reference | 27 | 27 | 27 | 27 | 27 | 55 | 27 | 27 | 27 | 27 | 27 | 27 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| phrases | **31** | **30** | 6 | **28** | 0 | **29** | 0 | - | 0 | 1 | 2 | 8 | 9 | 20 | 21 | 32 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\underbrace{\qquad\qquad\qquad}_{P}$ $\underbrace{\qquad\qquad\qquad}_{S}$

# Example



phrases | **31** | **30** | 6 | **28** | 0 | **29** | 0 | - | 0 | 1 | 2 | 8 | 9 | 20 | 21 | 32

$P$

$S$

mismatching symbols

31 30 28 29

zero-length factors

long factors

# Example

phrases

# Example



phrases

# Example



phrases

# Example

phrases

# Example



- predecessor structure to find the phrase that contains the start of the SA interval

# Example



- predecessor structure to find the phrase that contains the start of the SA interval

- $value[i] = \begin{cases} phrases[i], & \text{if } phraseLength[i] = 1 \\ reference[phrases[i]] + prevSaValue - n, & \text{otherwise} \end{cases}$

# Experiment

We compared our prototype (*rlzsa*) to other compressed indexes, replicating the experimental design used in the *r-index* paper (Gagie, Navarro, and Prezza, SODA 2018)

Datasets:

- boost — concatenated versions of GitHub's boost library — 600Mbyte

- DNA — concatenated copies of a DNA sequence of length 1000 with mutations — 600 Mbyte

- einstein — concatenated versions of Wikipedia's Einstein page — 600 Mbyte

- world — pdf files of CIA World Leaders from Jan 2003 to Dec 2009 — 45Mbyte

Search queries:

- 1000 patterns

- length = 8

# Experiment

We compared our prototype (*rlzsa*) to other compressed indexes, replicating the experimental design used in the *r-index* paper (Gagie, Navarro, and Prezza, SODA 2018)

Datasets:

- boost — concatenated versions of GitHub's boost library — 600Mbyte

- DNA — concatenated copies of a DNA sequence of length 1000 with mutations — 600 Mbyte

- einstein — concatenated versions of Wikipedia's Einstein page — 600 Mbyte

- world — pdf files of CIA World Leaders from Jan 2003 to Dec 2009 — 45Mbyte

Search queries:

- 1000 patterns

- length = 8

References:

- boost — 21 samples * 4096
- einstein — 2089 samples * 3072
- DNA — 11377 samples * 2048
- world — 498 samples * 4096

# Experimental results

# Comparison to r-index

|          | more space | much faster |
|----------|:----------:|:-----------:|
| boost    | 4.88       | 115.32      |
| DNA      | 13.19      | 149.10      |
| einstein | 4.65       | 96.27       |
| world    | 3.84       | 89.08       |

# Future work

- Reducing space:

  - Our prototype uses word-aligned units everywhere (16-, **32**-, 64-bit ints)

  - We can save space by using succinct representations instead (Elias-Fano for the predecessor structure, packed int vectors for phrases, etc.) (progress here)

  - Improved reference construction (we have some progress here already as well)

- Apply it to document (D) array

- Best of both worlds?

  - Is there a way to derive a hybrid of the **r-index** and **rlzsa**?

# Thank you!

r-index time $= O(occ * \log \log n)$

CDAWG time $= O(m(\log \log n + \log z) + pocc * \log^{\varepsilon} z + socc * \log \log n)$

**rlzsa** time $= O(\log \log n + occ + l_{max})$